

# Testing for Conditional Independence

- Outline:
- (i) Hardness of CI Testing
  - (ii) Kernel-based Testing (KCI/T)
  - (iii) Conditional Randomization Testing (CRT)
  - (iv) Classifier based Testing (CCIT)
  - (v) FOCI (Azadkia and Chatterjee)
  - (vi) Note on Tigranite (Runge et al).

i

## Hardness of CI Testing

Ref: The Hardness of CI Testing and the Generalized Covariance Measure, R. Shah and J. Peters, March 2021 (arXiv: 1804.07203v5).

Setup:  $(X, Y, Z)$  a triplet of random vectors s.t.

$$X \in \mathbb{R}^{d_x}, \quad Y \in \mathbb{R}^{d_y}, \quad Z \in \mathbb{R}^{d_z}, \quad \text{and}$$

are continuous r.v.s, i.e., their induced measure is absolutely continuous w.r.t. the Lebesgue measure.

$$\Sigma_0 = \left\{ \nu(\cdot, \cdot, \cdot) : \nu \text{ is a valid dist. over } (X, Y, Z), \right. \\ \left. \text{and is abs. cont. w.r.t. Lebesgue measure} \right\}$$

For any  $M > 0$ ,  $\mathcal{E}_{0,M} = \left\{ r(\cdot, \cdot, \cdot) : r \in \mathcal{E}_0 \text{ and } \begin{array}{l} \text{support}(X) \subseteq [-M, M] \\ \text{support}(Y) \subseteq [-M, M] \\ \text{support}(Z) \subseteq [-M, M] \end{array} \right\}$ .

cond. indep.  
continuous.

$\mathcal{P}_0 = \left\{ p(\cdot, \cdot, \cdot) : p \in \mathcal{E}_0 \text{ and } p \text{ s.t. } X \perp\!\!\!\perp Y \mid Z \right\}$

$\mathcal{Q}_0 = \mathcal{E}_0 \setminus \mathcal{P}_0$

$\mathcal{P}_{0,M} = \mathcal{P}_0 \cap \mathcal{E}_{0,M}$  ;  $\mathcal{Q}_{0,M} = \mathcal{Q}_0 \cap \mathcal{E}_{0,M}$ .

Composite Hypothesis Testing: (for any fixed  $M \in (0, \infty]$ )

(0) Null Hypothesis  $(X \perp\!\!\!\perp Y \mid Z)$  :  $\{p : p \in \mathcal{P}_{0,M}\}$

(1) Alternative  $(X \not\perp\!\!\!\perp Y \mid Z)$  :  $\{q : q \in \mathcal{Q}_{0,M}\}$ .

Goal: Given  $n$  iid samples  $\{(x_i, y_i, z_i)\}_{i=1}^n$ , we want to devise a test:

$\Psi_n : \mathbb{R}^{(d_x + d_y + d_z)n} \times [0, 1] \mapsto \{0, 1\}$ .

can be randomized, i.e.,  $\Psi_n(\{(x_i, y_i, z_i)\}_{i=1}^n, U)$   
independent Uniform  $[0, 1]$  r.v.  $\leftarrow$

→ Type I error

**Level of Test:** Test has valid level  $\alpha \in (0, 1)$  at sample size  $n$  if

$$\sup_{P \in \mathcal{P}_{0,M}} P_p(\Psi_n = 1) \leq \alpha$$

→ we would like  $\alpha \rightarrow 0$ .

$$(x_i, y_i, z_i) \sim P, \text{ iid } i=1, \dots, n$$

(i.e., this is the analog of  $P(\text{Test is wrong} | H_0)$  in binary hyp. testing.)

↖ null hypothesis.

→ 1 - (Type II error)

**Power of Test:** Test has power  $\beta \in (0, 1)$  at sample size  $n$  if

$$\inf_{Q \in \mathcal{Q}_{0,M}} P_q(\Psi_n = 1) \geq \beta.$$

→ we would like  $\beta \rightarrow 1$

$$(x_i, y_i, z_i) \sim q, \text{ iid } i=1, \dots, n$$

(This is the analog of  $P(\text{Test is correct} | H_1)$  in binary hyp. testing.)

The main hardness theorem below:

(Thm 2 in paper)

**Theorem 1:** Given any  $n \geq 1$ ,  $\alpha \in (0, 1)$ ,  $M \in (0, \infty]$  any any (randomized) test  $\Psi_n$  for which:

$$\sup_{P \in \mathcal{P}_{0, M}} \mathbb{P}_P(\Psi_n = 1) \leq \alpha$$

We have:  $\sup_{Q \in \mathcal{Q}_{0, M}} \mathbb{P}_Q(\Psi_n = 1) \leq \alpha.$

This says that there is no test for CI that can have low level and high power.

Total Variation (TV) Distance:  $P, Q \in \mathcal{E}_0$ . Then,

$$\|P - Q\|_{TV} = \sup_{A \in \mathcal{B}} \left| \mathbb{P}_P((X, Y, Z) \in A) - \mathbb{P}_Q((X, Y, Z) \in A) \right|$$

Borel  $\sigma$ -algebra over  $\mathbb{R}^{d_1 + d_2 + d_3}$

(Prop 5)

**Theorem 2:** For each  $M \in (0, \infty]$ ,  $\exists Q \in \mathcal{Q}_{0, M}$  st.

$$\inf_{P \in \mathcal{P}_{0, M}} \|P - Q\|_{TV} \geq 1/24.$$

Theorem 2 says that there is no  $p$  satisfying  $X|Y|Z$  that is close in TV distance to a specific  $q$  (construction below).

However, for any  $q$  (including the one constructed above), CI testing does not have any reasonable power!

Takeaway: TV distance is not capturing hardness of CI Testing. See discussion after Prop 5 in paper for more details.

Instead of going through the proof, I will present an example below that lies at the heart of their proof.

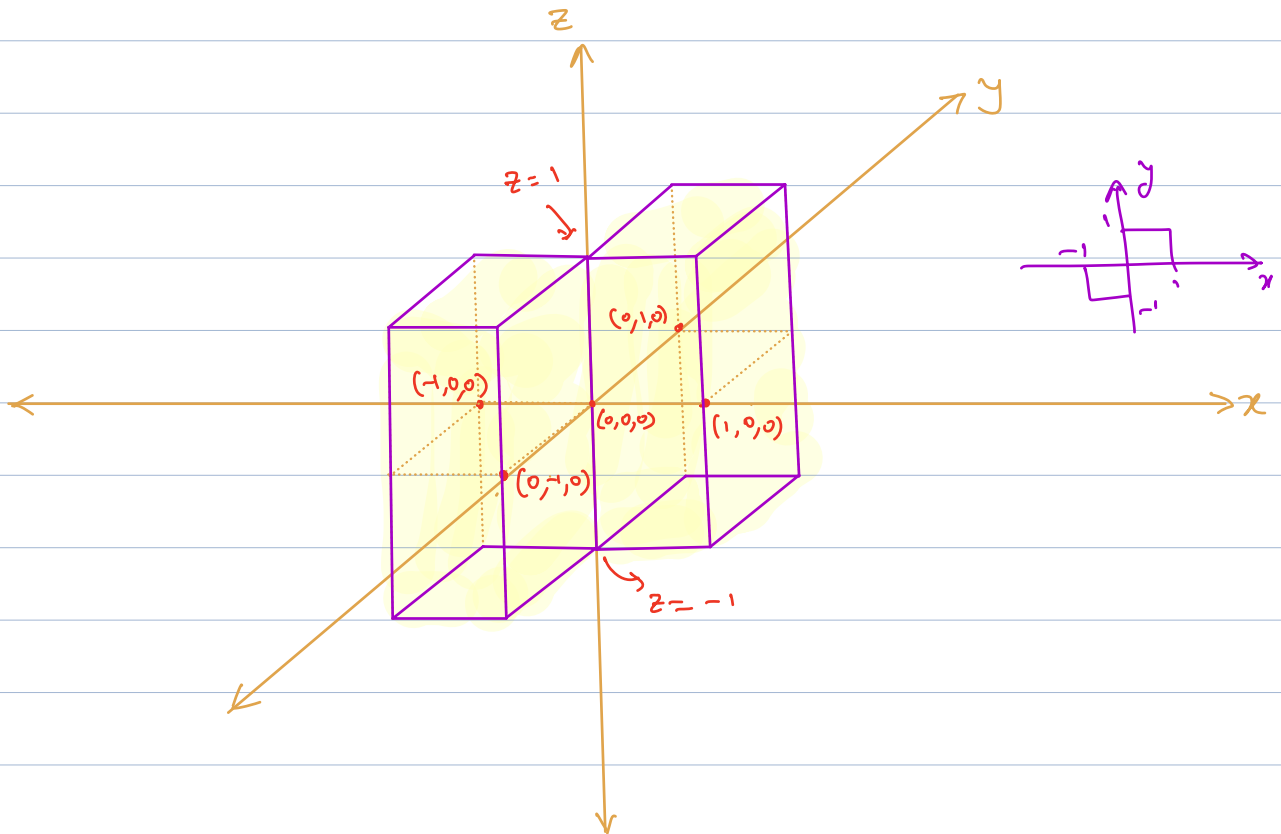
Construction for  $(X, Y, Z)$  scalar triplet over  $[-1, 1]^3$  to "show" Theorem 2 and plausibility of Theorem 1.

$X, Y, Z$  each over  $[-1, 1]$ , with their marginal dists:  $X \sim \text{Unif}[-1, 1]$ ,  $Y \sim \text{Unif}[-1, 1]$  and  $Z \sim \text{Unif}[-1, 1]$ .

$$q: (X \not\perp Y | Z)$$

$$Z \sim \text{Unif}[-1, 1]$$

$$(X, Y) \perp\!\!\!\perp Z \quad \text{with} \quad (X, Y) \sim \text{Uniform}([0, 1]^2 \cup [-1, 0]^2)$$



For the above  $q$ , observe that  $(X, Y)$  are correlated — If  $X > 0$ , then  $Y \sim \text{Unif}[0, 1]$  else if  $X < 0$ , then  $Y \sim \text{Unif}[-1, 0]$ .

$$\text{cov}(X, Y) = 1/4. \quad \Rightarrow \quad X \not\perp Y | Z$$

(note:  $X \not\perp Y$ , and  $(X, Y) \perp\!\!\!\perp Z$ )

Now, we construct  $\tilde{Z}$  a new r.v., with

$\tilde{Z} = f(x, Z)$  as follows: Consider the binary expansion of  $Z$ . In the  $m^{\text{th}}$  bit, delete the value, and replace it by '0' if  $x \geq 0$  and '1' if  $x < 0$ .

ex: Suppose the sample drawn  $z_i = 0.010010\dots$   
 (i.e.,  $z_i \approx \frac{1}{4} + \frac{1}{32} + \dots \approx 0.2812\dots$ ).

Then  $\tilde{z}_i = \begin{cases} 0.010010\dots & \text{if } x_i \geq 0 \\ 0.010011\dots & \text{if } x_i < 0. \end{cases}$

Now consider  $(X, Y, Z, \tilde{Z})$ .

Then, by choosing  $m$  to be large enough,  
 $|z_i - \tilde{z}_i| \ll 1$  (we can make the sample-wise distance to be arbitrarily small with error  $= \frac{1}{2^m}$ ).

∴ (a) From samples, it is very hard to distinguish between  $(X, Y, Z)$  and  $(X, Y, \tilde{Z})$ .

(b) The above construction encodes the sign of  $X$  noiselessly in  $\tilde{Z}$ . Furthermore, observe that there is no additional information in  $X$  about  $Y$ .

once we know it's sign, i.e.,

$$X \perp\!\!\!\perp Y \mid \tilde{Z}.$$

(c) Let  $p$  s.t.  $(X, Y, \tilde{Z}) \sim p$   
 $q$  s.t.  $(X, Y, Z) \sim q.$

The  $\|p - q\|_{TV} > \text{constant} > 0$   
 $\hookrightarrow$  independent of  $m.$

Why? Use a Borel set that focusses only on the  $m^{\text{th}}$  coordinate of the binary expansion of  $Z, \tilde{Z}$ , and some large enough set in  $X \times Y$  that includes both (I) and (II) quadrant. Then,  $Z$  and  $\tilde{Z}$  have a constant prob. of differing in the  $m^{\text{th}}$  component, and thus the TV distance is lower bounded  $> 0$ , independent of  $m.$

$\therefore$  Now making  $m$  large, we have fixed TV distance lower bound, but the samples get closer. Thus, it is hard to distinguish between  $p$  and  $q$ , even with many samples, so Thm 1 is plausible.



In fact, as shown in proof of Prop 2 in the paper, for any  $p \in \mathcal{P}_{0,1}$ ,  $\|p - q\|_{TV} > 1/24$ .

See Figure 1 in paper (Hardness of CI...) for idea of proof of Thm 1.

**Remark 1:** This crucially relies on  $Z$  being a continuous r.v., so that we can hide the sign of  $X$  with arbitrarily small perturbation in  $|z - z'|$ .

**Remark 2:** We saw that without assumptions, the CI composite hypothesis testing problem is untestable. The key issue is non-smoothness of  $p(x, y | z)$  in  $z$ , i.e., a small perturbation in  $z$  drastically changed the joint dist of  $(X, Y)$ .

In Theorem 2, we saw that we could have a large TV distance in  $\|p(x, y, z) - q(x, y, z)\|_{TV}$  and still have issues, so controlling the distance between  $p \in \mathcal{P}_{0,M}$  and  $q \in \mathcal{Q}_{0,M}$  does not seem

useful. Indeed as remarked in Neykov, Balakrishnan and Wasserman 2021, we need to control smoothness of the conditional  $p(x, y|z)$ . A similar condition was imposed for the CCIT algorithm to have guarantees in Sen, Suresh, Karthikeyan, Dimakis, Shakkottai 2017.

The condn would look something like:

$$\|q(x, y|z) - q(x, y|z')\|_{TV} < L \|z - z'\|_2 \quad \left. \vphantom{\|q(x, y|z) - q(x, y|z')\|_{TV}} \right\} \begin{array}{l} \text{alternate} \\ \text{hypothesis.} \\ (x \perp\!\!\!\perp y|z) \end{array}$$

$$\|P(x|z) - P(x|z')\|_{TV} < L \|z - z'\|_2 \quad \left. \vphantom{\|P(x|z) - P(x|z')\|_{TV}} \right\} \begin{array}{l} \text{null} \\ \text{hypothesis.} \\ (x \perp\!\!\!\perp y|z) \end{array}$$

$$\|P(y|z) - P(y|z')\|_{TV} < L \|z - z'\|_2$$

Ref:

Minimax optimal conditional independence testing, M. Neykov, S. Balakrishnan and L. Wasserman, arXiv:2001.03039, 2021.

Model power conditional independence test, R. Sen, A.T. Suresh, K. Shanmugan, A. Dimakis and S. Shakkottai, NeurIPS 2017.

ii

## Kernel-based CI Testing

Ref: Partial associative measures and an application to qualitative regression, J. Daudin, Biometrika 1980

A kernel statistical test for independence, A. Gretton, K. Fukumizu, C. Teo, L. Song, B. Scholkopf, A. Smola, NeurIPS 2007.

Kernel-based conditional independence test and application in

causal recovery, K. Zhang, J. Peters, D. Janzing and B. Schölkopf, UAI 2011. (arXiv 1202.3775)

### Warm up: Jointly Gaussian $(X, Y, Z)$

For jointly Gaussian r.v.s, CI is equivalent to independence testing of residuals, i.e.,

$$X \perp\!\!\!\perp Y \mid Z \iff R_{XZ} \perp\!\!\!\perp R_{YZ}$$

$$\begin{array}{c} \text{partial correlation} \\ \text{coefficient} \end{array} \longleftarrow \rho_{X,Y|Z} = 0.$$

$$R_{XZ} = X - E[X|Z], \quad R_{YZ} = Y - E[Y|Z].$$

### Partial correlation coefficient

For jointly Gaussian distributions, conditional independence can be tested through regression, using partial correlation coefficient.

Recall that for Gaussians, it suffices to check second order statistics for independence, i.e.,

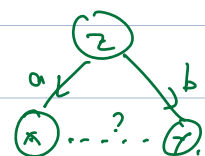
$$X \perp Y \iff \rho_{XY} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} = 0.$$

The partial correlation coefficient generalizes this to allow tests for CI for joint Gaussian r.v.s.

$$X \perp Y \mid Z \iff \rho_{XY.Z} = 0.$$

We will define  $\rho_{XY.Z}$  below. WLOG, all r.v.s below have zero mean ( $E[X] = E[Y] = E[Z] = 0$ )

$$\text{let } a = \underset{\alpha}{\text{argmin}} E[(X - \alpha Z)^2]$$



$$b = \underset{\beta}{\text{argmin}} E[(Y - \beta Z)^2]$$

i.e., regress  $X$  and  $Y$  each separately on  $Z$ .

$$\hat{X} = E[X|Z] = aZ, \quad \hat{Y} = E[Y|Z] = bZ$$

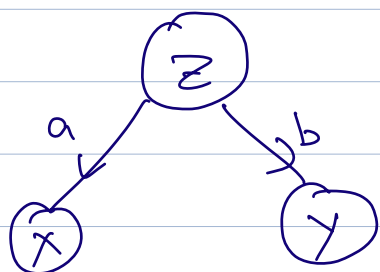
are the associated least-squares estimators for  $X, Y$  respectively.

$R_{xz} = (x - E[x|z])$        $R_{yz} = (y - E[y|z])$   
 are the associated residuals.

Then,  $\rho_{xy.z} = \rho_{R_x R_y} = \frac{\text{Cov}(R_{xz}, R_{yz})}{\sqrt{\text{Var}(R_{xz}) \text{Var}(R_{yz})}}$

Then,  $\rho_{xy.z} = 0 \iff x \perp\!\!\!\perp y \mid z$ .

To see why:



$$z = N_1,$$

$$x = a z + N_2$$

$$y = b z + N_3$$

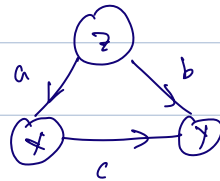
$N_1, N_2, N_3 \sim N(0, 1)$   
 mutually independent

In the model above, it is clear from the discussions in Note 2 (d-separation) that  $x \perp\!\!\!\perp y \mid z$ .

Now, let us compute  $\rho_{xy.z}$ . Here  $\hat{x} = a z$ ,  $\hat{y} = b z$ ,  $R_{xz} = (x - a z)$ ,  $R_{yz} = (y - b z)$ .

$$E[R_{xz}R_{yz}] = ab - ab - ab - ab = 0.$$

To check the other way, consider a model with  $c \neq 0$ , and check that  $\rho_{x \cdot z} \neq 0$



More generally, if  $(Z_1, Z_2, \dots, Z_m)$  are the covariates (i.e., possible confounding r.v.s), we regress  $X$  on  $(Z_1, \dots, Z_m)$  and also regress  $Y$  on  $(Z_1, \dots, Z_m)$ . Then compute  $\rho_{x \cdot z}$  be the correlation coefficient of the residuals.

Note: If the r.v.s are not jointly Gaussian, then we cannot relate partial correlation to independence in either direction. (see example 7.9 in text)

Beyond Joint Gaussians: A general criterion for CI

(Dawid '80) A useful characterization of CI is the following:  $X \in \mathbb{R}^{d_1}$ ,  $Y \in \mathbb{R}^{d_2}$ ,  $Z \in \mathbb{R}^{d_3}$

$$F: \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \mapsto \mathbb{R}, \quad g: \mathbb{R}^{d_1} \times \mathbb{R}^{d_3} \mapsto \mathbb{R}$$

$$\text{s.t. } E[F(x, z)^2] < \infty, \quad E[g(y, z)^2] < \infty$$

$$\tilde{F}(x, z) = F(x, z) - E[F(x, z) | z]$$

$$\tilde{g}(y, z) = g(y, z) - E[g(y, z) | z]$$

$$X \perp\!\!\!\perp Y | Z \iff E[\tilde{F}(x, z) \tilde{g}(y, z)] = 0$$

$$\forall F, g \text{ square-integrable s.t. } E[F(x, z) | z] = 0$$

$$\text{and } E[g(y, z) | z] = 0 \text{ a.s.}$$

Other equivalent criteria: (Daudin 80, Zhang et al 11)

$$\mathcal{E}_{XZ} = \left\{ \tilde{F} \text{ sq. int. : } E[\tilde{F}(x, z) | z] = 0 \text{ a.s.} \right\}$$

$$\mathcal{E}_{YZ} = \left\{ \tilde{g} \text{ sq. int. : } E[\tilde{g}(y, z) | z] = 0 \text{ a.s.} \right\}$$

$$\mathcal{E}'_{YZ} = \left\{ \tilde{g}' : \tilde{g}'(y, z) = g'(y) - E[g'(y) | z], \right. \\ \left. g' \text{ sq. int. function of } y \right\}$$

$$(i) \quad X \perp\!\!\!\perp Y | Z$$



$$(ii) \quad E[\tilde{f}(x, z) \tilde{g}(y, z)] = 0 \quad \forall \quad \begin{matrix} \tilde{f} \in \mathcal{E}_{xz}, \\ \tilde{g} \in \mathcal{E}_{yz} \end{matrix}$$

$$(iii) \quad E[\tilde{f}(x, z) g(y, z)] = 0 \quad \forall \quad \begin{matrix} \tilde{f} \in \mathcal{E}_{xz}, \\ g \text{ sq. int. i.e.,} \\ E[g(y, z)^2] < \infty \end{matrix}$$

$$(iv) \quad E[\tilde{f}(x, z) \tilde{g}'(y, z)] = 0 \quad \forall \quad \begin{matrix} \tilde{f} \in \mathcal{E}_{xz}, \\ \tilde{g}' \in \mathcal{E}'_{yz} \end{matrix}$$

$$(v) \quad E[\tilde{f}(x, z) g'(y)] = 0 \quad \forall \quad \begin{matrix} \tilde{f} \in \mathcal{E}_{xz}, \\ g' \text{ r.t. } E[g'(y)^2] < \infty \end{matrix}$$

Constructing a Test: **KCIT** (Operationalizing (iv) above)

$$X \in \mathcal{X} = \mathbb{R}^{d_x}, \quad Y \in \mathcal{Y} \in \mathbb{R}^{d_y}, \quad Z \in \mathcal{Z} \in \mathbb{R}^{d_z}.$$

$k_x$  positive definite kernel over RKHS  $\mathcal{H}_x$ , analogous  $k_y, k_z$ .

**Characteristic kernel**:  $k_x$  is characteristic if  $E_{x \sim p, q}[f(x)] = E_{x \sim p}[f(x)] \quad \forall f \in \mathcal{H}_x \implies p = q$

Implication is that we can use a kernel to test for one of the CI criteria above (generalizes residue tests).

The KCIT approach (Zhang et al. 2011) is the following:

Samples:  $x = (x_1, \dots, x_n)$ ,  $y = (y_1, \dots, y_n)$ ,  $z = (z_1, \dots, z_n)$

$$\tilde{K}_x = \left( I - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right) K_x \left( I - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right)$$

centered  
sample kernel

$$\begin{pmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & K_x(x_i, x_j) & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{pmatrix}$$

$$= V_x \Lambda_x V_x^T$$

non-negative eigenvalues

$$\Psi_x = \begin{pmatrix} \Psi_{x,1} & \Psi_{x,2} & \dots & \Psi_{x,n} \end{pmatrix}, \quad \text{where } \Psi_{x,i} = \sqrt{\lambda_{x,i}} V_{x,i}$$

Similarly for  $y, z$  and  $\tilde{x} = (x, z)$

fix  $\varepsilon > 0$  (regularizer parameter).

$$R_z = \varepsilon (\tilde{K}_z + \varepsilon I)^{-1}, \quad \tilde{K}_{z|z} = R_z \tilde{K}_z R_z$$

$$\tilde{K}_{y|z} = R_z \tilde{K}_y R_z$$

## Constructing Test Statistic with finite samples, with asymptotic characterization

$$T_{CI}^{(n)} = \frac{1}{n} \text{trace}(\tilde{K}_{x|z} \cdot \tilde{K}_{y|z})$$

**Theorem (Prop 5 in KCI $\pi$ ):** Under null hypothesis ( $x \perp\!\!\!\perp y | z$ ),  $T_{CI}^{(n)}$  has the same asymptotic dist. (i.e. conv. in distribution) as  $\tilde{T}_{CI}^{(n)}$ , where

$$\tilde{T}_{CI}^{(n)} = \frac{1}{n} \sum_{k=1}^{n^2} \tilde{\lambda}_k r_k^2, \text{ where}$$

$\tilde{\lambda}_k$  are eigenvalues of  $\tilde{w} \tilde{w}^T$ , and  $\tilde{w}$  defined by

$$\tilde{w} = \begin{pmatrix} \psi_{x|z} \\ \psi_{y|z} \end{pmatrix}$$

vector defined analogous as above for  $\psi_x$ , but for the appropriate centered kernel matrices.

and  $\{\sigma_k\}$  are iid  $N(0,1)$  r.v. (or equivalently,  $r_k^2$  are iid  $\chi^2$ -distributed r.v.s).

$$\begin{matrix} \nearrow P(T \geq t | H_0) \\ < P < \alpha \end{matrix}$$

**Punchline:** We can operationalize the above for p-value testing. Simulate  $\tilde{T}_{CI}^{(n)}$  on generate samples. Compute test statistic  $T_{CI}^{(n)}$  from data. See where  $T_{CI}^{(n)}$  falls in  $\tilde{T}_{CI}^{(n)}$  to compute significance level to reject null hypothesis.

Remark (and caution): All the properties and associated tests are based on population (i.e., exact distribution) based characterization of conditional independence. Unlike the "usual" proofs where finite sample versions follow from "standard" concentrations, CI is quite different.

A "generic" way to break a CI test based on exact distributions/population statistics is the following: Recall the example above, where the sign of  $X$  was embedded in the  $m^{\text{th}}$  bit, thus slightly perturbing  $Z$  to  $\tilde{Z}$  through a sample-path coupling.

$H_0: \{P_m\}: X \perp\!\!\!\perp Y \mid Z$	$P_m: \text{encode sign of } X \text{ in } m^{\text{th}} \text{ bit}$
$H_1: q: Z \perp\!\!\!\perp (X, Y)$ with $X \not\perp\!\!\!\perp Y \mid Z$	$X \perp\!\!\!\perp Y \mid \tilde{Z}_m$

But  $q$  and  $P_m$  are very close in samples, i.e.,  $(X, Y, Z) \sim q$  and  $(X, Y, \tilde{Z}_m) \sim P_m$  are only  $1/2^m$  apart in  $|Z - \tilde{Z}_m|$ .

Goal: Given a tester with a CI characterization that is based on infinite samples/exact-dist., we want to construct

a hard counter-example that will break it for any finite  $n$  samples.

Given:  $n$  samples, CI tester.

Adversary: Picks one of  $q$  or  $P_n$  (i.e., sign of  $X$  in  $n^{\text{th}}$  bit).

Problem: Impossible to reliably distinguish between these with  $n$  samples, as signal strength ( $1/2^n$ ) is buried in sampling noise. Thus, for any fixed  $n$ , tester fails.

Work around: Restrict null and alternatives to have Lipschitz continuous conditionals. Then, all the algorithms that we discuss based on asymptotic properties will "work".

(iii)

Conditional Randomization Test (CRT)

Ref: Panning for gold: Model-X knockoffs for high-dimensional controlled variable selection, E. Candès, Y. Fan, L. Janson and J. Lv, Journal of the Royal Statistical Society, Series B, 2018 (arXiv: 1610.02351).

(Motivation is feature selection in regression).

Setting:  $(X_1, X_2, \dots, X_d, Y)$

Goal: Want to check if:  $X_1 \perp\!\!\!\perp Y \mid (X_2, \dots, X_d)$ .

Assume: We have exact knowledge of the joint dist. of  $(X_1, \dots, X_d)$ , i.e.,  $p(x_1, x_2, \dots, x_d)$ , and we can generate samples from this dist.

Suppose that  $T = g(X_1, X_2, \dots, X_d, Y)$  is a test statistic

(e.g.:  $T$  is the (regularized) regression coefficient of  $X_1$ ,  
when we regress  $Y$  on  $(X_1, \dots, X_d)$  .  
e.g. lasso

Auxillary Variable: Let  $X_1^* \sim p(\cdot \mid x_2, \dots, x_d)$  be a conditionally indep. r.v., i.e., given  $X_2 = x_2, X_3 = x_3, \dots, X_d = x_d$ , we generate a sample from the dist.  $p(z \mid x_2, \dots, x_d)$ . With this construction, observe that:  
↓  
dummy variable.

$$X_1^* \mid (X_2, \dots, X_d, Y) \stackrel{d}{=} X_1 \mid (X_2, \dots, X_d)$$

The idea is that if  $X_1 \perp\!\!\!\perp Y \mid (X_2, X_3, \dots, X_d)$ , then  $T = g(X_1, X_2, \dots, X_d, Y)$  and  $T^* = g(X_1^*, X_2, \dots, X_d, Y)$  would have identical dist. for every value of  $(X_2, X_3, \dots, X_d, Y)$ . This is summarized in their lemma (Lemma 4.1 in Candès et al. 2018) below:

**Key Observation** (Lemma 4.1 in paper): Under the null hypothesis,

$$T \mid (X_2, X_3, \dots, X_d, Y) \stackrel{d}{=} T^* \mid (X_2, X_3, \dots, X_d, Y)$$

This provides an asymptotic test for CI: (a) Simulate a large number of samples  $X_1^*$ , as above; (b) Compute the test statistic for each; (c) Compute test statistic for the real  $X_1$ ; (d) Compute p-value.

Remark 1: This test is based on a population property under the null hypothesis, namely the original and simulated test statistics are identically distributed.

However, when dealing with finite samples, some issue that was discussed re. untestability arises. We need to control smoothness of conditionals to provide any finite sample guarantees (note that the paper

does not address finite sample guarantees).

Remark 2: Candès et al 2018 provide an alternate approach based on knockoff- $X$  random variables, that is computationally much lighter (in the CRT, we need to generate the test statistic for each sample). The power with knockoffs is somewhat worse; please read their paper for details.

## Classifier-based Conditional Independence Test

Ref: Model-powered Conditional Independence Test, R. Sen, A. T. Suresh, K. Sharmugam, A. Dimakis, S. Shakkottai, NeurIPS 2017. (arXiv:1709.06138).

Setting:  $(X, Y, Z)$ ,  $X \in \mathbb{R}^{d_x}$ ,  $Y \in \mathbb{R}^{d_y}$ ,  $Z \in \mathbb{R}^{d_z}$

$$H_0: X \perp\!\!\!\perp Y | Z$$

$$H_1: X \not\perp Y | Z.$$

Given:  $3n$  iid samples

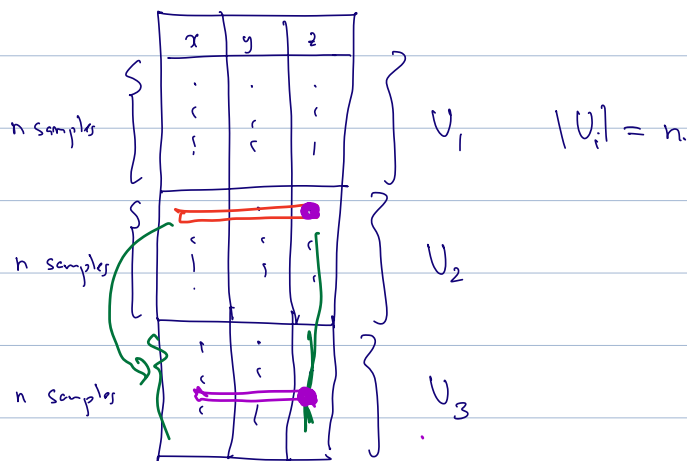
$$\{(x_i, y_i, z_i)\}_{i=1}^{3n} \sim p(x, y, z)$$

Under  $H_0$ :  $p(x, y, z) = p^{CI}(x, y, z) \triangleq p(x|z) p(y|z) p(z)$



Assumptions:  $Z$  a continuous r.v. that satisfies smoother, in both  $p(z)$  (marginal) and conditional  $P(y|z)$ .  
 (characterized through max eigenvalue of fisher information matrix of  $y$  w.r.t  $z$ )

Algorithm:



① Partition into 3 groups of  $n$ -samples each.

② Nearest neighbour bootstrap:

a. For each  $u = (x, y, z) \in U_2$ , find  $(x', y', z') \in U_3$  s.t.  $z' \in 1\text{-NN}(z)$

↙  
nearest neighbor  
w.r.t.  $l_2$ -norm.

b. Construct  $u' = (x, y', z)$  from above step.

(Note that if  $z = z'$ , then  $u' \sim p(x|z)p(y|z)p(z)$  irrespective of whether  $H_0$  is true or  $H_1$  is true.)

Let  $\phi(x, y, z)$  s.t.  $U' \sim \phi(x, y, z)$ .

$$p^{CI}(x, y, z) = p(x|z) p(y|z) p(z).$$

Note  $\phi \approx p^{CI}$   $\left( \begin{array}{l} \text{Theorem 1 in Sen et. al.} \\ \text{characterizes TV distance between} \\ \phi \text{ and } p^{CI} \text{ as roughly } \left( \frac{1}{n^{1/2d_z}} \right). \end{array} \right)$

So far:  $\star$   $n$  iid samples  $(x, y, z) \sim p$  (labeled  $U_1$ )

$\star$   $n$  samples  $(x, y, z) \sim \phi$  (labeled  $U_2'$ )

Note:

(i) Samples in  $U_2'$  are not iid, because of the 1-NN process, that dips into the same pool  $U_2$  for each  $(x, y, z) \in U_2$  to construct  $(x, y', z)$ .

(ii) Samples in  $U_2'$  are only approximately CI, i.e.,  $\phi(x, y, z) \approx p(x|z) p(y|z) p(z)$ .

③ Train a classifier: Dataset:  $U_1 \sim p(x, y, z)$   
 $U_2' \sim \phi(x, y, z)$

① label samples in  $U_1$  as '1' (n samples)  
Samples in  $U_2'$  as '0' (n nearly-CI samples)

Intuition:  $u \in U_1$ ,  $u' \in U_2'$  are almost identically distributed under  $\mathcal{H}_0 (X \perp\!\!\!\perp Y | Z)$ . However  $u \in U_1$  and  $u' \in U_2'$  have different dist. under  $\mathcal{H}_1$  (thus, note that this property is again a population property and thus we need smoothness & conditionals to exploit this property.)

② Split  $D = (\{U_1: 1\} \cup \{U_2': 0\})$  into train and test ( $D_{tr}$  and  $D_{test}$ ), with  $|D_{tr}| = |D_{test}| = n$ .

Using  $D_{tr}$ , train a binary classifier, s.t. the classification function class is rich enough (formally, the risk of the best classifier from this class is close to that under the Bayes optimal classifier, see Thm 2 in Sen et. al.).

③ Using this classifier, evaluate  $D_{test}$  using

ERM, with risk function:  $R_n = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\{g(u) \neq l\}}$   
test sample  $\rightarrow$  classifier  $\rightarrow$  label

Result (Implication of Thm 2 in Sen et. al.):

Under  $H_0 (x \perp\!\!\!\perp y | z)$ :  $R_n \approx 0.5$

with high enough prob.

Under  $H_1 (x \not\perp\!\!\!\perp y | z)$ :  $R_n \leq 0.5\delta + o(n)$

$d_{TV}(P, P^{(z)}) \geq 1 - \delta$